

Anomaly Attribution with Likelihood Compensation

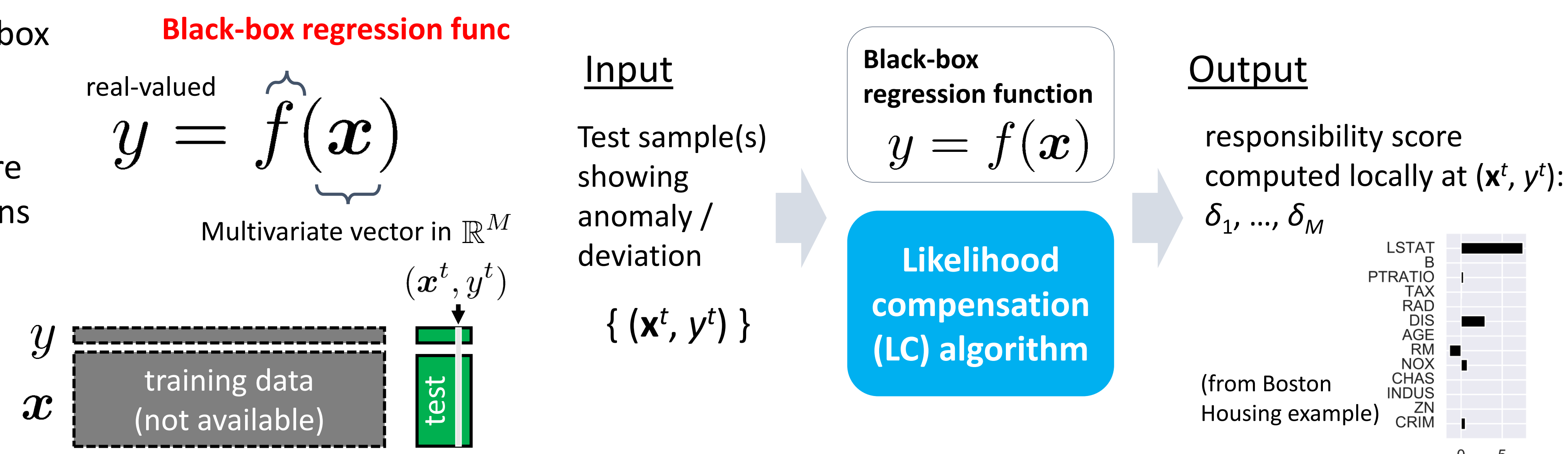
Tsuyoshi (“Ide-san”) Idé, Amit Dhurandhar, Jiří Navrátil, Moninder Singh, Naoki Abe

IBM T. J. Watson Research Center

Slides available: <https://ide-research.net>

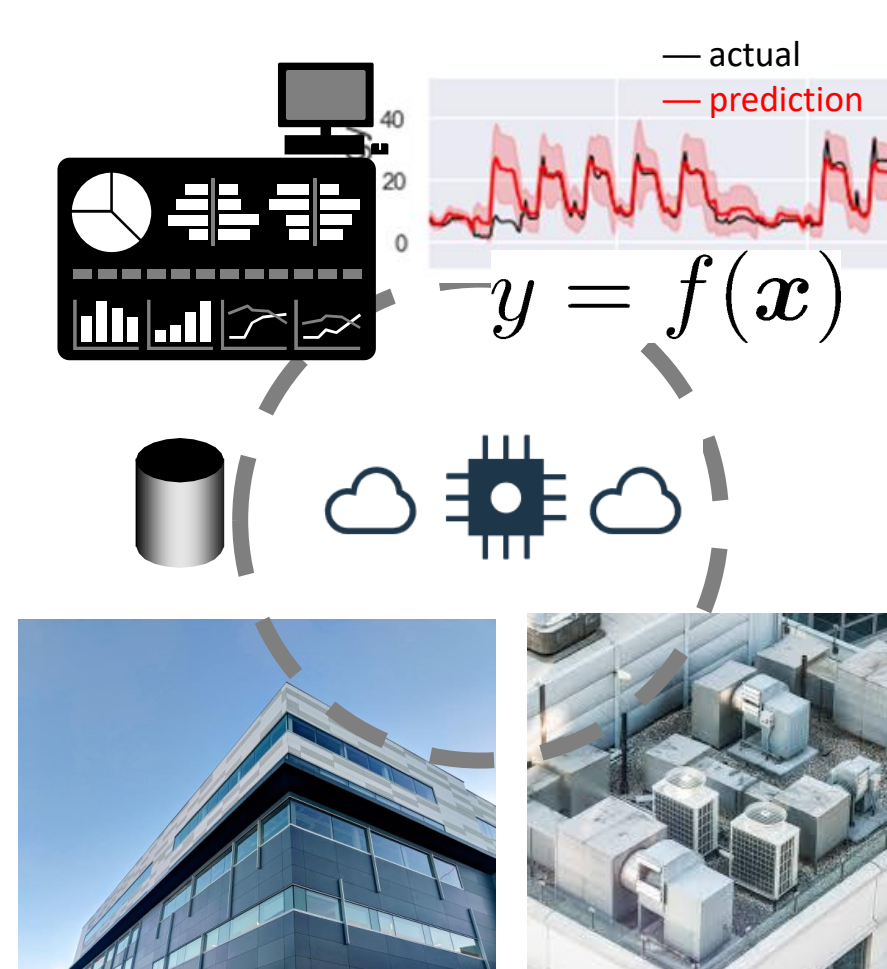
Technical task: anomaly attribution for black-box regression

- Task:** Attribute deviation from black-box prediction $f(\mathbf{x})$ to each input variable
- Background:** Most of XAI methods are designed to explain $f(\mathbf{x})$, not deviations
- Solution:** New notion of “likelihood compensation”

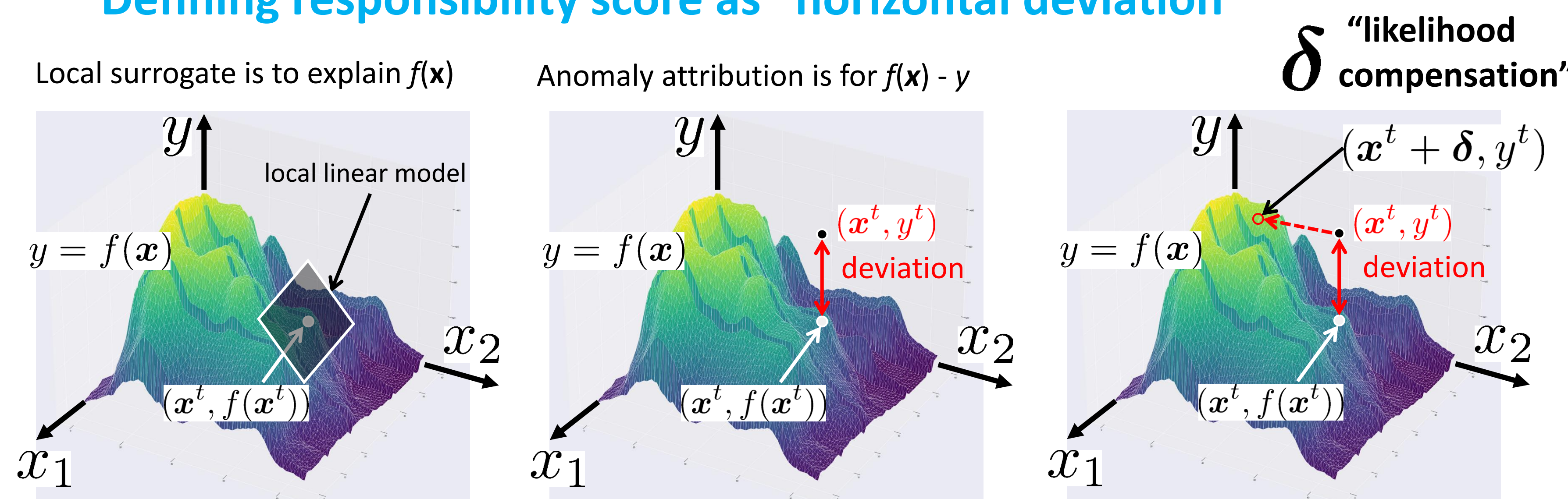


Use-case example: Building energy management

- y : building energy consumption
- \mathbf{x} : Temperature, humidity, day of week, month, room occupancy, etc.
- Why black-box?
 - Multiple players (AI vendor / Sier / HVAC constructor)
 - Proprietary technologies



Local surrogate model for $f(\mathbf{x})$ alone cannot explain deviations. Defining responsibility score as “horizontal deviation”



LC as optimal perturbation to \mathbf{x}^t

- Likelihood compensation δ :** A perturbation to \mathbf{x}^t such that $\mathbf{x}^t + \delta$ achieves the best possible fit to the model
- LC seeks a best possible fit by correcting \mathbf{x}^t under a certain regularization

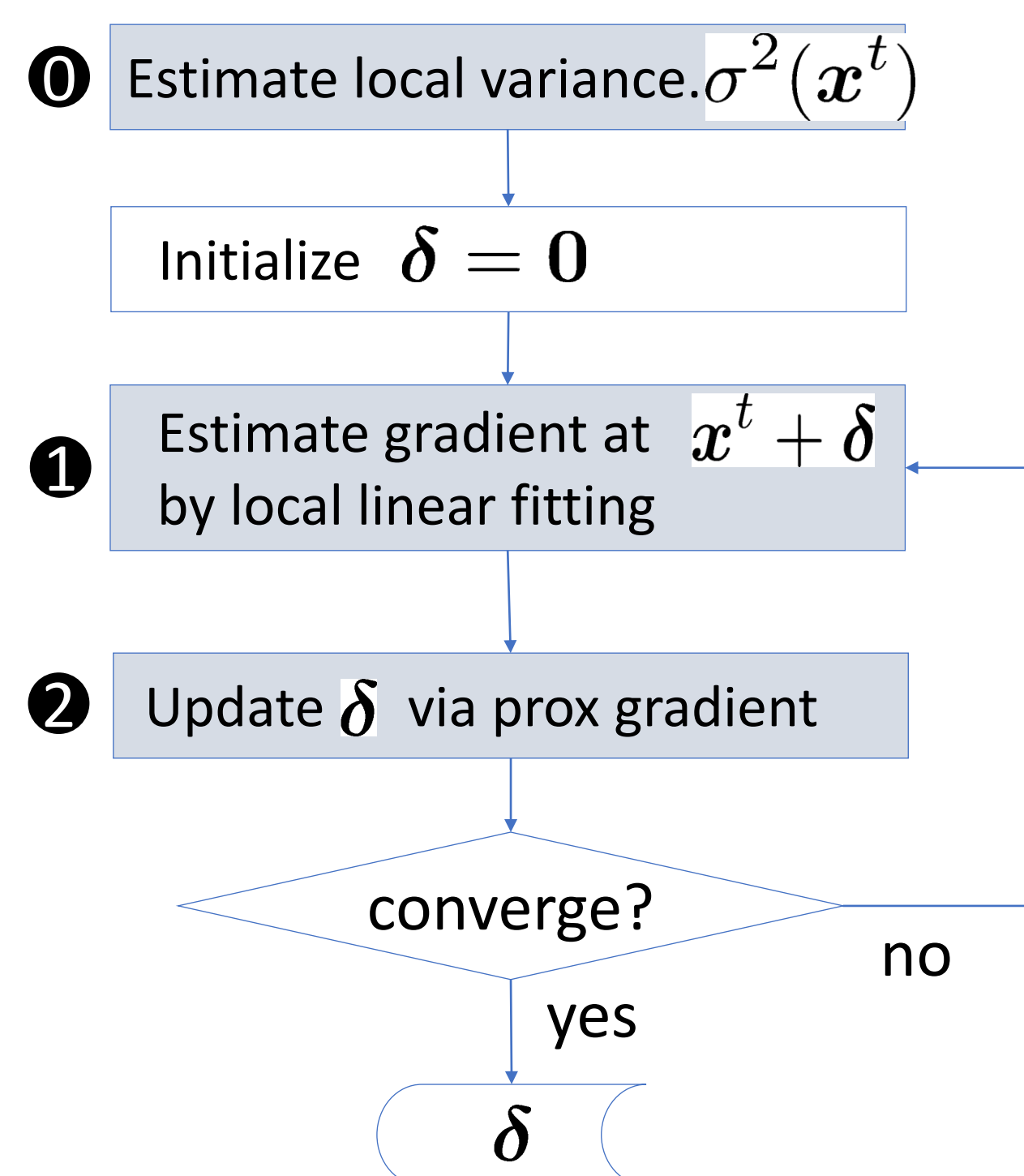
$$\delta = \arg \max_{\delta} \left[\ln \left\{ p(y^t | f(\mathbf{x}^t + \delta)) p(\delta) \right\} \right],$$

Gaussian elastic net

- The main optimization problem

$$\min_{\delta} \left\{ \frac{1}{N_{\text{test}}} \sum_{t=1}^{N_{\text{test}}} \frac{[y^t - f(\mathbf{x}^t + \delta)]^2}{2\sigma_t^2} + \frac{1}{2} \lambda \|\delta\|_2^2 + \nu \|\delta\|_1 \right\},$$

(λ and ν are constant)



0. Local variance estimation

- No extra test sample available \rightarrow use constant variance (no choice!)

$$p(y | \mathbf{x}) = \mathcal{N}(y | f(\mathbf{x}), \sigma^2)$$

- Some amount of test samples available \rightarrow locally weighted maximum likelihood

$$\sigma^2(\mathbf{x}^t) = \max_{\sigma^2} \sum_{n=1}^{N_{\text{heldout}}} w_n(\mathbf{x}^t) \left\{ \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{(y^{(n)} - f(\mathbf{x}^{(n)}))^2}{2\sigma^2} \right\},$$

Gaussian kernel defined for the specific test sample \mathbf{x}^t

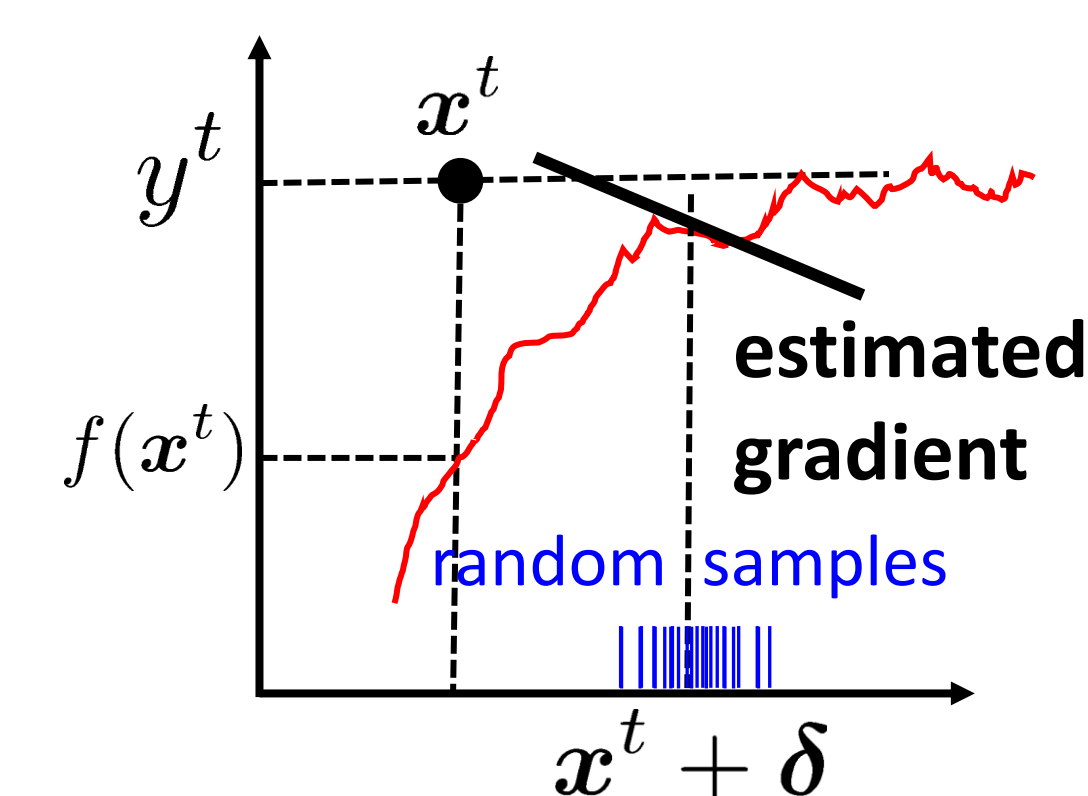
1. Local gradient estimation of f

- $f(\mathbf{x})$ is black-box but we need the gradient:

$$\frac{1}{N_{\text{test}}} \sum_{t=1}^{N_{\text{test}}} \frac{y^t - f(\mathbf{x}^t + \delta)}{\sigma_t^2} \left\| \frac{\partial f(\mathbf{x}^t + \delta)}{\partial \delta} \right\| + \lambda \delta$$

- Gradient $\left\| \frac{\partial f(\mathbf{x}^t + \delta)}{\partial \delta} \right\|$

\rightarrow Local sampling & linear fit



2. Proximal gradient update for δ

- With the numerically estimated gradient, the problem now looks like L_1 -regularized convex-ish optimization

- Updating rule from δ^{old} using prox gradient-like algorithm

$$J(\delta) \triangleq \frac{1}{N_{\text{test}}} \sum_{t=1}^{N_{\text{test}}} \frac{[y^t - f(\mathbf{x}^t + \delta)]^2}{2\sigma_t^2} + \frac{1}{2} \lambda \|\delta\|_2^2$$

$$\delta = \arg \min_{\delta} \left\{ J(\delta^{\text{old}}) + (\delta - \delta^{\text{old}}) \langle \nabla J(\delta^{\text{old}}) \rangle + \frac{1}{2\kappa} \|\delta - \delta^{\text{old}}\|_2^2 + \nu \|\delta\|_1 \right\}$$

$$= \text{PROX}_{\kappa\nu\|\cdot\|_1} \left(\delta^{\text{old}} - \kappa \langle \nabla J(\delta^{\text{old}}) \rangle \right)$$

Has an analytic solution! (\rightarrow paper)

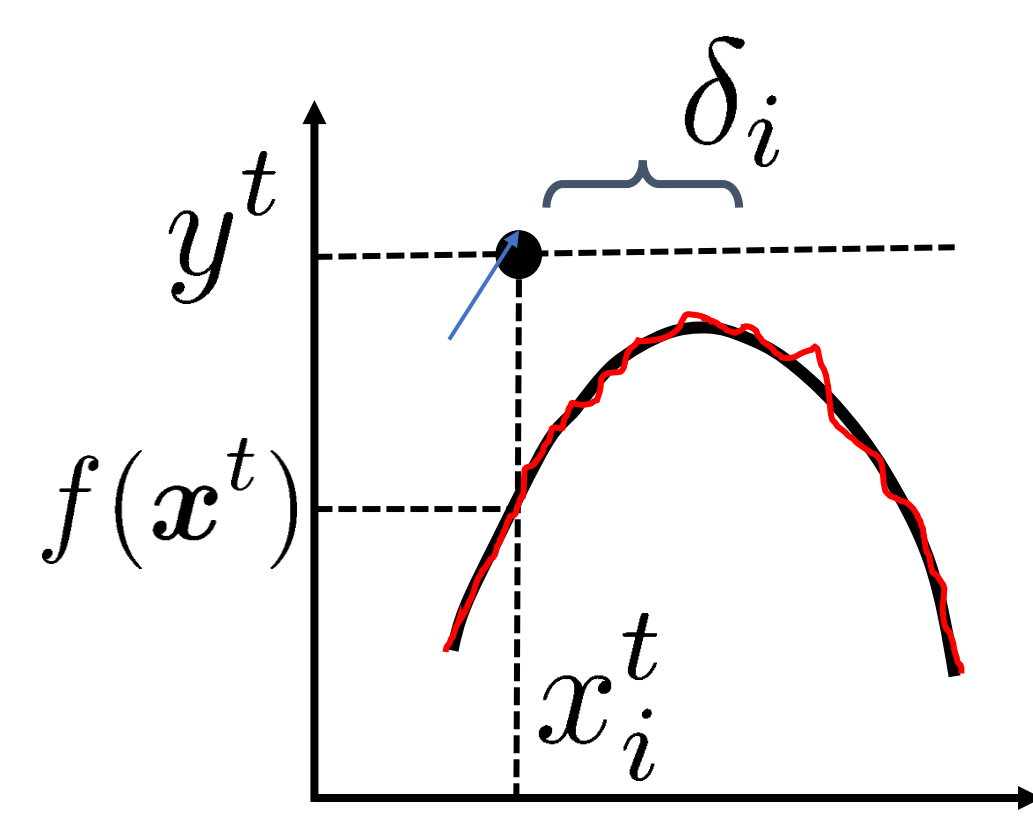
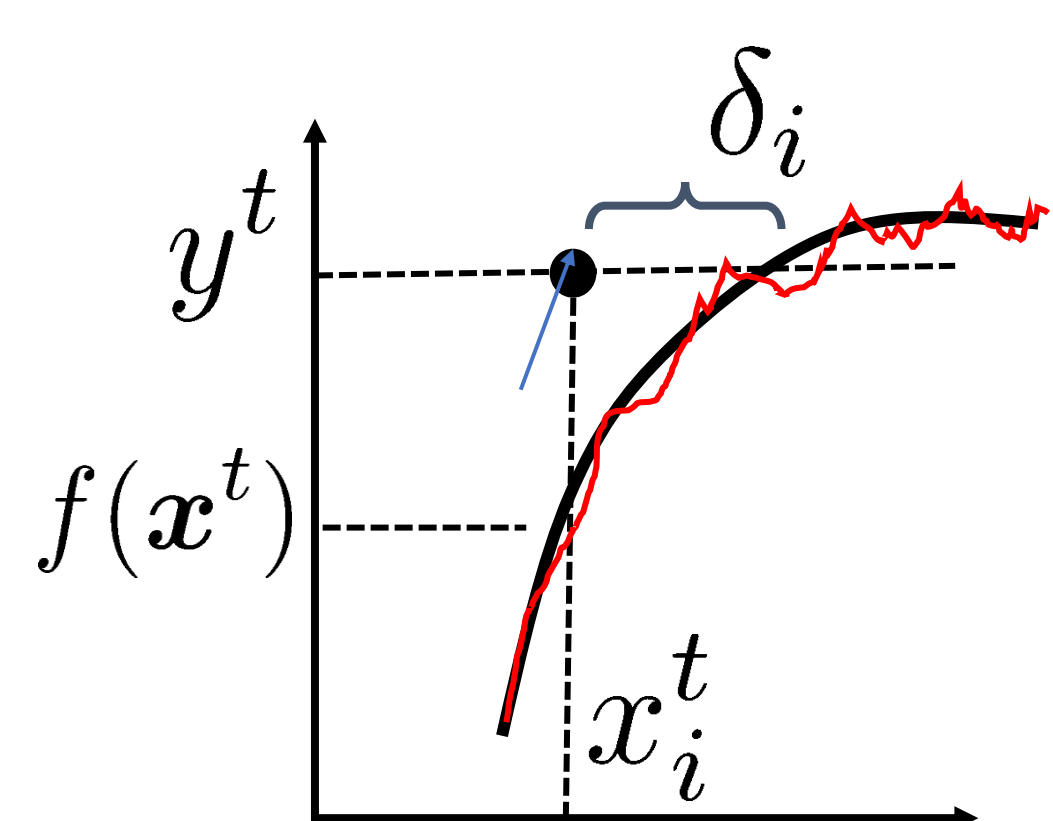
Condition of convergence – where the intuition of “horizontal deviation” comes from

- The prox gradient-like update converges when

$$\frac{1}{N_{\text{test}}} \sum_{t=1}^{N_{\text{test}}} \frac{y^t - f(\mathbf{x}^t + \delta)}{\sigma_t^2} \left\| \frac{\partial f(\mathbf{x}^t + \delta)}{\partial \delta} \right\| \approx 0$$

Cond. 1: |deviation| = 0

Cond. 2: |gradient| = 0



Existing methods for anomaly attribution

- Few applicable approaches
- Existing methods are mostly for explaining $f(\mathbf{x})$ in classification, often in white-box setting

- Possible baselines:

1: Z-score

$$a_i(\mathbf{x}^t) = \frac{x_i^t - \text{mean}_i}{\text{stddev}_i}$$

2: LIME [Ribeiro 18], extended to explain $f(\mathbf{x}^t) - y^t$ rather than $f(\mathbf{x}^t)$

3: Shapley value [Strumbelj+ 14], extended to explain $f(\mathbf{x}^t) - y^t$

Comparison in the building energy use-case

- One month-worth building energy data
 - y : energy consumption
 - \mathbf{x} : time of day, temperature, humidity, sunrad, day of week (one-hot encoded)
- Top: Overall anomaly score
 - \equiv energy deviation
- Bottom 3: Responsibility score of each input variable
 - Only LC captures meaningful patterns

